

# Alibaba Cloud E-MapReduce

Quick Start

Issue: 20181119

# Legal disclaimer

---

Alibaba Cloud reminds you to carefully read and fully understand the terms and conditions of this legal disclaimer before you read or use this document. If you have read or used this document, it shall be deemed as your total acceptance of this legal disclaimer.

1. You shall download and obtain this document from the Alibaba Cloud website or other Alibaba Cloud-authorized channels, and use this document for your own legal business activities only. The content of this document is considered confidential information of Alibaba Cloud. You shall strictly abide by the confidentiality obligations. No part of this document shall be disclosed or provided to any third party for use without the prior written consent of Alibaba Cloud.
2. No part of this document shall be excerpted, translated, reproduced, transmitted, or disseminated by any organization, company, or individual in any form or by any means without the prior written consent of Alibaba Cloud.
3. The content of this document may be changed due to product version upgrades, adjustments, or other reasons. Alibaba Cloud reserves the right to modify the content of this document without notice and the updated versions of this document will be occasionally released through Alibaba Cloud-authorized channels. You shall pay attention to the version changes of this document as they occur and download and obtain the most up-to-date version of this document from Alibaba Cloud-authorized channels.
4. This document serves only as a reference guide for your use of Alibaba Cloud products and services. Alibaba Cloud provides the document in the context that Alibaba Cloud products and services are provided on an "as is", "with all faults" and "as available" basis. Alibaba Cloud makes every effort to provide relevant operational guidance based on existing technologies. However, Alibaba Cloud hereby makes a clear statement that it in no way guarantees the accuracy, integrity, applicability, and reliability of the content of this document, either explicitly or implicitly. Alibaba Cloud shall not bear any liability for any errors or financial losses incurred by any organizations, companies, or individuals arising from their download, use, or trust in this document. Alibaba Cloud shall not, under any circumstances, bear responsibility for any indirect, consequential, exemplary, incidental, special, or punitive damages, including lost profits arising from the use or trust in this document, even if Alibaba Cloud has been notified of the possibility of such a loss.
5. By law, all the content of the Alibaba Cloud website, including but not limited to works, products, images, archives, information, materials, website architecture, website graphic layout, and webpage design, are intellectual property of Alibaba Cloud and/or its affiliates. This intellectual property includes, but is not limited to, trademark rights, patent rights, copyrights, and trade

secrets. No part of the Alibaba Cloud website, product programs, or content shall be used, modified, reproduced, publicly transmitted, changed, disseminated, distributed, or published without the prior written consent of Alibaba Cloud and/or its affiliates. The names owned by Alibaba Cloud shall not be used, published, or reproduced for marketing, advertising, promotion, or other purposes without the prior written consent of Alibaba Cloud. The names owned by Alibaba Cloud include, but are not limited to, "Alibaba Cloud", "Aliyun", "HiChina", and other brands of Alibaba Cloud and/or its affiliates, which appear separately or in combination, as well as the auxiliary signs and patterns of the preceding brands, or anything similar to the company names, trade names, trademarks, product or service names, domain names, patterns, logos, marks, signs, or special descriptions that third parties identify as Alibaba Cloud and/or its affiliates).

6. Please contact Alibaba Cloud directly if you discover any errors in this document.



# Generic conventions

Table -1: Style conventions

Style	Description	Example
	This warning information indicates a situation that will cause major system changes, faults, physical injuries, and other adverse results.	 <b>Danger:</b> Resetting will result in the loss of user configuration data.
	This warning information indicates a situation that may cause major system changes, faults, physical injuries, and other adverse results.	 <b>Warning:</b> Restarting will cause business interruption. About 10 minutes are required to restore business.
	This indicates warning information, supplementary instructions, and other content that the user must understand.	 <b>Note:</b> Take the necessary precautions to save exported data containing sensitive information.
	This indicates supplemental instructions, best practices, tips, and other content that is good to know for the user.	 <b>Note:</b> You can use <b>Ctrl + A</b> to select all files.
>	Multi-level menu cascade.	<b>Settings &gt; Network &gt; Set network type</b>
<b>Bold</b>	It is used for buttons, menus, page names, and other UI elements.	Click <b>OK</b> .
Courier font	It is used for commands.	Run the <code>cd /d C:/windows</code> command to enter the Windows system folder.
<i>Italics</i>	It is used for parameters and variables.	<code>bae log list --instanceid Instance_ID</code>
[] or [a b]	It indicates that it is a optional value, and only one item can be selected.	<code>ipconfig [-all -t]</code>
{ } or {a b}	It indicates that it is a required value, and only one item can be selected.	<code>swich {stand   slave}</code>

# Contents

---

<b>Legal disclaimer</b> .....	<b>I</b>
<b>Generic conventions</b> .....	<b>I</b>
<b>1 Prerequisites</b> .....	<b>1</b>
<b>2 Create E-MapReduce</b> .....	<b>3</b>
2.1 Quick start.....	3
2.2 Create a cluster.....	4
2.3 Create a job.....	8
2.4 Create an execution plan.....	12
<b>3 Cluster specifications</b> .....	<b>15</b>

# 1 Prerequisites

---

Before creating E-MapReduce, you need to complete the following prerequisites:

1. Apply for an Alibaba Cloud account.

Before applying for an E-MapReduce cluster, you need to have an Alibaba Cloud account to identify yourself through the entire Alibaba Cloud ecosystem. This account can be used not only for E-MapReduce clusters, but also to activate Alibaba Cloud services, such as [Object Storage Service \(OSS\)](#) and [ApsaraDB for RDS \(RDS\)](#).

If you do not have any Alibaba Cloud account, see [Register Cloud Account](#).

2. Create an AccessKey (optional).

You must create at least one AccessKey according to the following steps:

- a. Log on to the [Alibaba Cloud console](#).
- b. Click **AccessKeys**.

**Note:**

If a security prompt dialog box appears, click **Continue to manage AccessKey**.

## Security Tips

---



AccessKey of your cloud account is the secret key to access Alibaba Cloud APIs. Since the AccessKey has full permissions of your cloud account, please make sure you keep it well. To avoid the AccessKey being used by others to cause [Sensitive information leakage](#), do not release your AccessKey to any external channels (for example, Github). We strongly recommend you use the AccessKeys of RAM users in API calls, according to [Alibaba Cloud account security best practices](#).

Continue to manage AccessKey

Get Started with Sub Users's AccessKey

- c. Click **Create AccessKey**.
  - d. AccessKey is created successfully.
3. Activate Alibaba Cloud OSS.

E-MapReduce will store your job logs and run logs in the Alibaba Cloud OSS storage space, so you need to [Sign up for OSS](#). And create a bucket in the same area where you expect to create the cluster, see [Create a bucket](#).

#### 4. Enable high-end models (optional)

If you want to use models with 8 cores or more in clusters charged by the Pay-As-You-Go billing method, apply for opening it in ECS first. [Apply for high-end models](#).

## 2 Create E-MapReduce

---

### 2.1 Quick start

In this tutorial, you will learn how clusters, jobs, and execution plans are used in E-MapReduce. You will also learn how to create a Spark Pi job, run it successfully in the cluster, and check the results.

**Note:**

Make sure you have completed the necessary [prerequisites](#).

1. Create a cluster.
  - a. At the top of the [EMR console](#), click **Cluster Management > Create Cluster**.
  - b. Software configurations.
    - A. Use the latest EMR product version, such as **EMR-3.13.0**.
    - B. Use default software configurations.
  - c. Hardware configurations.
    - A. Select **Pay-as-You-Go**.
    - B. If there is no security group, enter a name to create new security group.
    - C. Select **4 vCPU 8G** for the **Master Instance Type**.
    - D. Select **4 vCPU 8G** for the **Core Instance Type** (one instance).
    - E. Keep others in the default status.
  - d. Basic configurations.
    - A. Enter the name of the cluster.
    - B. Select a log path to save job logs and make sure that the **Running Logs** button is turned on. In the cluster related region, [create an OSS bucket](#).
    - C. Enter the password.
    - D. Click **Next**.
  - e. Click **Create** to create a cluster.
2. Create a job.
  - a. At the top of the page, click the **Data Platform** tab.
  - b. In the upper right corner of the **Projects** page, click **New Project**.
  - c. In the **New Project** panel, enter the project name and project description, and click **Create**.

- d. At the right side of the corresponding project, click **Design Workflow**.
- e. In the **New Job** dialog box, enter the job name, job description, and select Spark as the job type.

Once the job type is selected, it cannot be modified.

- f. Click **OK**.

**Note:**

You can also create subfolders, rename folders, and delete folders by right-clicking on the folders.

- g. In the **Content** box, enter parameters as follows.

```
--class org.apache.spark.examples.SparkPi --master yarn-client
--driver-memory 512m --num-executors 1 --executor-memory 1g
--executor-cores 2 /usr/lib/spark-current/examples/jars/spark-
examples_2.11-2.1.1.jar 10
```

**Note:**

The `/usr/lib/spark-current/examples/jars/spark-examples_2.11-2.1.1.jar` jar file name is decided by the Spark version in the cluster, for example, if Spark version is 2.1.1, it should be `spark-examples_2.11-2.1.1.jar`, if Spark version is 2.2.0, then file name is `spark-examples_2.11-2.2.0.jar`.

- h. Click **Run**.

3. View job logs and confirm the results.

After the job runs, at the bottom of the page, in the **Running Records** tab, you can view running logs of the job. Click **Logs** to jump to the detailed log page of the job, you can see information in **Job Instance**, **Runtime Log**, and **YARN Containers**.

## 2.2 Create a cluster

In this tutorial, you will learn how to create a cluster.

### Enter the cluster creation page

1. Log on to the [Alibaba Cloud E-MapReduce console](#).
2. Complete RAM authorization. For procedure, see [Role authorization](#).
3. Select a region for a cluster. The region cannot be changed once the cluster is created.
4. Click **Create Cluster** to create a cluster.

## Cluster creation process

To create a cluster, follow the below steps:

### 1. Step one: Software configuration

Configuration description

- **EMR version:** Select the latest version by default.
- **Cluster type:** Currently E-MapReduce provides the following cluster types:
  - Hadoop clusters, provide semi-managed ecosystem components:
    - Hadoop, Hive, and Spark that offline store and compute distributed data at scale.
    - SparkStreaming, Flink, and Storm that are stream processing systems.
    - Presto and Impala, for running interactive analytics queries.
    - Oozie and Pig.
  - Druid clusters, provide semi-managed, real-time interactive analysis services, query large amount of data in millisecond latency, and support for multiple data intake methods . Used with services such as EMR Hadoop, EMR Spark, OSS, and RDS, Druid clusters offer real-time query solutions.
  - Data Science clusters, are mainly for big data and AI scenarios, providing Hive and Spark offline big data, and TensorFlow model training.
  - Kafka clusters, are taken as a semi-managed distributed message system of high throughput and high scalability, providing a complete service monitoring system that can keep a stable running environment.
- **Include configurations:** Use the default configuration. You can add, start and stop services in the management interface later.
- **High security mode:** In this mode, you can set the Kerberos authentication of the cluster. This feature is unnecessary for clusters used by individual users. It is turned off by default.
- **Enable custom setting:** You can specify a JSON file to change software configuration before you start a cluster.

### 2. Hardware configuration

Configuration description

- Billing configuration
  - **Billing method:** Pay-As-You-Go is used in testing scenarios. Subscription production clusters can be created after all tests are verified.

- Network configuration
  - **Zone:** Generally use the default zone.
  - **Network type:** By default, the Virtual Private Cloud (VPC) network is selected which requires you to enter a VPC and a VSwitch. If you haven't created a network, go to the [VPC console](#) to create them.
  - **VPC:** Select the region of the VPC network.
  - **VSwitch:** Select a zone for VSwitch under the corresponding VPC. If no VSwitch is available in this zone, then you must create a new one.
  - **Security group name:** Generally, no security group exists when you create a cluster for the first time. Enter a name to create a new security group. If you already have a security group in use, you can choose to use it directly here.
- Cluster configuration
  - **High availability:** When enabled, two master instances in the Hadoop cluster are used to ensure the availability of the resource manager and name node. HBase clusters support high availability by default. When enabled, a master instance is used to ensure high availability.
  - Master node
    - **Master instance type:** Select an instance type as required. For information about instance types, see [Instance type families](#).
    - **System disk type:** Select a disk as required.
    - **System disk size:** We recommend that you select 120G at least.
    - **Data disk type:** Select a disk as required.
    - **Data disk size:** We recommend that you select 80G at least.
    - **Master instances:** It is set 1 master instance by default.
  - Core node
    - **Core instance type:** Select an instance type as required. For information about instance types, see [Instance type families](#).
    - **System disk type:** Select a disk as required.
    - **System disk size:** We recommend that you select 80G at least.
    - **Data disk type:** Select a disk as required.
    - **Data disk size:** We recommend that you select 80G at least.
    - **Core instances:** It is set 2 instances by default. You can adjust as required.

- Task instance group: It is turned off by default.

### 3. Basic configuration

#### Configuration description

- Basic information

- **Cluster name:** The cluster name can contain Chinese characters, English letters (uppercase and lowercase), numbers, hyphens (-), and underscores (\_), with a length limit between 1-64 characters.

- Running logs

- **Running logs:** The function for saving running logs is turned on by default. In the default state, you can select the OSS directory location to save running logs. You must activate OSS before using this function. Cost depends on the number of uploaded files. We recommend that you open the OSS log saving function, which helps in debugging and error screening.

- **Log path:** OSS path for saving the log.

- **Uniform Meta Database:** We recommend you disable this feature for the moment.

- Permission settings: Use default settings.

- Logon settings

- **Remote logon:** It is turned on by default to enable security group port 22.

- **Logon password:** Set the logon password at the master node. The logon password must contain English letters (both uppercase and lowercase letters), numbers, and special characters (!@#%\$%^&\*) with a length limit between 8-30 characters.

### Purchase list and cluster cost

Confirm the configured items and the billing in the configuration list.

### Confirm creation

After all configurations are set, the **Create** button is highlighted. Verify the information, and click **Create** to create clusters.



#### Note:

- If it is a Pay-As-You-Go cluster, the cluster is created immediately, and you are taken back to the **Overview** page where you can see a cluster in **Initializing** status. It takes several minutes to create the cluster. After creation, the cluster is switched to the **Idle** status.

- For subscribed clusters, the cluster is not created until the order is generated and paid.

### Creation failure

If cluster creation failed, the message **Cluster creation failed** appears on the cluster list page. The reason for the failure can be seen when the pointer is placed on the red exclamation point.

No handling is required because the corresponding computing resources are not created. The cluster is automatically hidden after three days.

## 2.3 Create a job

In this tutorial, you will learn how jobs are created in E-MapReduce.

To run a computing task, you need to define a job first according to the following steps:

1. Log on to [Alibaba Cloud E-MapReduce console](#).
2. Select the region where the job is created.
3. At the top of the page, click **Old EMR Scheduling**.
4. In the upper right corner of the page, click **Create Job**.

Create job✕

---

**\* Name :**

Length: 1 to 64 characters. Only Chinese characters, English letters, numbers '-', and '\_' are allowed

**\* Type :**

<input checked="" type="radio"/> Spark	<input type="radio"/> Hadoop	<input type="radio"/> Hive	<input type="radio"/> Pig
<input type="radio"/> Sqoop	<input type="radio"/> Spark SQL	<input type="radio"/> Shell	

**\* Parameter :**

+ Select OSS pathoss console Upload

**\* Actual execution :** **spark-submit**

Fail retry  Yes  No

**\* Failure policy :**

Pause current execution plan

Continue execution of next job

OK

Cancel

5. Enter the job **Name**.

6. Select a job **Type**.

7. Enter **Parameters** of the job. Parameters must include full information of the jar package used by the job, data input and output addresses of the job, and some command line parameters, that is, all your parameters in the command line must be completed in this field. You can click **Select OSS path** to select an OSS resource path. For parameter configurations of all job types, see the Job chapter in *Cite LeftUser GuideCite Right*.

- 8. Actual execution:** The actual executed command for the job on ECS will be displayed. If you copy the displayed command, the command can be run directly in the command line environment of the E-MapReduce cluster.
- 9. Fail retry:** If you select **Yes**, you can set the number of retries and each retry interval. It is set **No** by default.
- 10.Failure policy:** Pausing the current execution plan will pause the entire execution plan after this job fails and will wait for your handling. Continuing to execute the next job will ignore this error and continue to execute the next job after this job fails.
- 11.** Click **OK** to complete the creation.

### Job Example

This is a Spark job, where relevant parameters, input and output paths, are set in application parameters.

**Note:**

This example is only for reference.

## Create job



\* Name :   
 Length: 1 to 64 characters. Only Chinese characters, English letters, numbers '-', and '\_' are allowed

\* Type :  Spark  Hadoop  Hive  Pig  
 Sqoop  Spark SQL  Shell

\* Parameter : 

```
spark-submit --class org.apache.spark.examples.SparkPi --master
yarn-client --driver-memory 512m --num-executors 1 --executor-
memory 1g --executor-cores 2 /opt/apps/spark-1.6.1-bin-
hadoop2.7/lib/spark-examples-1.6.1-hadoop2.7.2.jar 100
```

  
[+ Select OSS path](#)

\* Actual execution : **spark-submit** spark-submit --class  
 org.apache.spark.examples.SparkPi --master yarn-client --driver-  
 memory 512m --num-executors 1 --executor-memory 1g --  
 executor-cores 2 /opt/apps/spark-1.6.1-bin-hadoop2.7/lib/spark-  
 examples-1.6.1-hadoop2.7.2.jar 100

\* Failure policy :  Pause current execution plan  
 Continue execution of next job

## OSS and ossref

The **oss://** prefix indicates that the data path is directed to an OSS path, which specifies the operation path similar to **hdfs://** when reading/writing the data.

The **ossref://** is also directed to an OSS path, however, it will be used to download the corresponding code resource to the local disk, and then replace the path in the command line with the local path. It is intended for easily running some native codes without the need to log on to the computer to upload the code and the dependent resource package.

The *ossref://xxxxxx/xxx.jar* parameter in the preceding example represents the jar of job resources. This jar is stored in OSS. When the system is running, the jar will be downloaded to clusters for running automatically in E-MapReduce. The two *oss://xxxx* and another two values behind a jar file appearing as parameters are passed to the main class in a jar file.

**Note:**

The **ossref** cannot be used to download excessive data resources, otherwise it will lead to failure of the cluster job.

## 2.4 Create an execution plan

In this tutorial, you will learn how execution plans are created in E-MapReduce.

After job creation, if you want to run the job defined on the cluster, you need to create an execution plan. An execution plan can contain more than one job, and you can define their order. For example, if one of your scenarios is: prepare data > process data > clean up data, you can define three jobs named **prepare-data**, **process-data**, and **cleanup-data**, and then create an execution plan to include these three jobs.

The steps to create an execution plan are as follows:

1. Log on to the [Alibaba Cloud E-MapReduce console](#).
2. Select a region for your cluster.
3. At the top of the navigation bar, click **Old EMR Scheduling**.
4. In the right-side panel, click **Execution plan**.
5. In the upper-right corner, click **Create an execution plan**.
6. In the **Select the cluster mode** pane, there are two options,
  - **Create as needed**: The cluster is created on demand, and the cluster is released immediately after the execution of the plan. If this option is selected, you must execute the same steps as creating a cluster, configure a cluster as needed, and then click OK.
  - **Existing clusters**: Select a cluster from the cluster you have created. If this option is selected, you will go to the **Select the cluster** panel where you can set an associated cluster and see the cluster information.

Create an execution plan
✕

---

1 : Select the
2 : Select the
3 : Configure the
4 : Configure the

\* Associated cluster : HBase ▼

\* Cluster info :

Status : Idle

Type : Hadoop

Running time : 5hour(s)36minute(s)43second(s)

created time : 2018/10/24 14:57:13

Prev
Next
Cancel

Execution plans can only be submitted to clusters in **Running** and **Idle** status.

#### 7. Click **Next** to enter the **Configure the job** page.

The page will display the previously defined job list, and the job list to be run as per the newly created execution plan on the right. Select jobs on the left to populate the right as per the execution order. You can click the question mark to view the detailed parameters.

Create an execution plan
✕

---

1 : Select the cluster mode
2 : Select the cluster
3 : Configure the job
4 : Configure the scheduling

Job list

Search

ID	Name	Type
J-224FB0D8 84AE6DEF	test_wordcount (?)	Hadoop
J-6C48C6EB 9D5AF93A	Pi (?)	Spark
J-B21F6DBC 51A68F7F	wcr (?)	Hadoop

Previous page
Next page

Current page : rank1page , total3records ,

Configured job

ID	Name	Type

Prev
Next
Cancel

#### 8. Click **Next**.

9. In the **Configure the scheduling** pane, enter the execution plan name and select a scheduling policy.

- **Manually executed:** Execution plans are executed manually.
- **Periodic scheduling execution plan:** Set the **scheduling cycle** and **first execute time** to run execution plans automatically.

Create an execution plan
✕

---

1 : Select the

2 : Select the

3 : Configure the

4 : Configure the

\* Name :

\* Scheduling policy : Periodic scheduling execution plan ▼

\* Set the scheduling cycle : day(s) ▼

\* Set the scheduling cycle : per 1 day(s)

\* first execute time : 2018-10-25 09 : 47

First run time 2018-10-25 9:47

Subsequent intervals 1 day(s) run 1 Times

When the execution plan is completed, you can make the relevant settings for the execution plan in the execution plan management page.

Prev
OK
Cancel

10. Click **OK**.

## 3 Cluster specifications

---

The first step to use E-MapReduce is to configure an appropriate Hadoop cluster. To select a E-MapReduce configuration type, you must take into account the big data scenarios of the enterprise to estimate Inoata amount, the reliability requirement of services, and corporate budgets.

### Big data scenarios

At present, E-MapReduce product primarily meets the following requirements of big data scenarios for enterprises:

- Batch processing features high disk throughput, high network throughput, and low real-time requirements. If there is large amounts of data to be processed, but the requirement for real-time processing is not high, you can select services such as MapReduce, Pig, and Spark. For scenarios with low memory requirements, consider the demand for CPU and memory for large jobs and the network demand of shuffle to select appropriate types.
- Ad hoc queries: data scientists or data analysts use ad hoc query tools to retrieve data. For scenarios with features of high real-time query capacity, high disk throughput, and high network throughput, use EMR Impala and Presto services. For scenarios with high memory requirements, consider data volume and the number of concurrent queries.
- For stream computing scenarios with high network throughput and compute-intensive requirements, use E-MapReduce Flink, Spark Streaming, and Storm services.
- For scenarios of message queues, high disk throughput, high network throughput, large amounts of memory consumption, and the scenarios where the storage doesn't depend on HDFS, use E-MapReduce Kafka. In order to avoid the impact on Hadoop, E-MapReduce will separate Kafka and Hadoop into two clusters.
- For data cold backup scenarios, the requirements for computing and disk throughput are not high, but low costs of data cold backup are required. We recommend that you use E-MapReduce D1 instance to perform data cold backup. The instance storage cost in the local D1 disk is USD 0.003/month/GB.

### E-MapReduce node

E-MapReduce has three *instance types*: Master, Core, and Task.

E-MapReduce storage can use ultra disks, SSD disks, and local disks. Disk performance comparison is as follows: SSD disks > local disks > ultra disks.

The E-MapReduce underlying storage supports OSS (standard OSS only) and HDFS. The data availability of OSS is higher than that of HDFS. The data availability of OSS is 99.99999999%, while that of HDFS is 99.99999%.

The storage price is estimated roughly as follows:

- The cost of local disk instance storage is USD 0.003/GB/month.
- The cost of OSS standard storage is USD 0.02/GB/month.
- The cost of ultra disks is USD 0.05/GB/month.
- The cost of SSD disks is USD 0.143/GB/month.

### Type selection of E-MapReduce

- Type selection of master nodes

Master nodes are used to deploy the master processes of Hadoop, such as NameNode and ResourceManager.

We recommend that you enable HA on the production cluster. HA is available for components such as E-MapReduce HDFS, YARN, Hive, and HBase. We recommend that you enable **High Availability** in the **Hardware configuration** step when you create the production cluster. If **High Availability** is not enabled when you purchase the product, it cannot be enabled later.

The master node is primarily used to store HDFS metadata and component log files. It is compute-intensive without high requirement for disk IO. HDFS metadata is stored in memory, and it is recommended that the memory should be 16 GB or above based on the number of files.

- Type selection of core nodes

The core node is primarily used to store data, perform calculations, and run DataNode and Nodemanager.

If the data volume of HDFS (3 backups) is greater than 60 TB, we recommend that you use local disk instance (ECS.D1, ECS.D1NE). The disk capacity is  $(\text{number of CPU cores}/2) \times 5.5\text{TB} \times \text{number of instances}$ . For example, if you purchase four 8-core D1 instances, the disk capacity is  $8/2 \times 5.5 \times 4 = 88$  TB. Because HDFS uses 3 backups, you can buy at least 3 local disk instances. Considering data reliability and disk damage, we recommend that you purchase at least 4 instances.

If the data volume of HDFS is less than 60TB, you can consider ultra disks and SSD disks.

- Type selection of task nodes

Task nodes are mainly used to supplement the CPU and memory, which can enhance the computing capability of Core nodes. The nodes do not store data or run DataNode. You can estimate the number of instances based on CPU and memory requirements.

### **E-MapReduce lifecycle**

E-MapReduce supports auto-scaling, which can quickly [expand a cluster](#). It can flexibly adjust the configuration of cluster nodes or [upgrade or downgrade instance configurations](#).

### **Available zones**

To ensure efficiency, it's better to deploy EMR and the business system in the same region and zone. For more information, see [regions and zones](#).